

Sécurité Appliquée, PVP

Devoir maison à rendre avant le 11/12 à 23h45

Jean-François COUCHOT
couchot [arobase] femto-st [point] fr

1 Partie pratique

Le TD2 doit se terminer sur l'idée que pour des valeurs entières, le mécanisme géométrique est plus utile que le mécanisme laplacien car la variance du bruit est toujours plus faible dans le premier cas que dans le second. On rappelle que :

— $\text{Var}[V] = 2 \left(\frac{\Delta}{\epsilon}\right)^2$ si V suit une loi de Laplace($0, \frac{\Delta}{\epsilon}$)

— $\text{Var}[G] = 2 \frac{e^{-\frac{\epsilon}{\Delta}}}{\left(1 - e^{-\frac{\epsilon}{\Delta}}\right)^2}$ si G est la variable aléatoire réelle (de bruit géométrique) définie comme au TD2.

L'objectif ici est d'illustrer ceci.

1.1 Comparaison (théorique) des variances

On considère ici une sensibilité Δ égale à 1.

Dans une première figure, afficher les deux courbes de variances pour ϵ qui varie de 0.1 à 10. On prendra une échelle logarithmique pour le budget ϵ . On devrait obtenir une figure semblable à celle donnée à la figure 1. On donnera le code python qui génère cette figure.

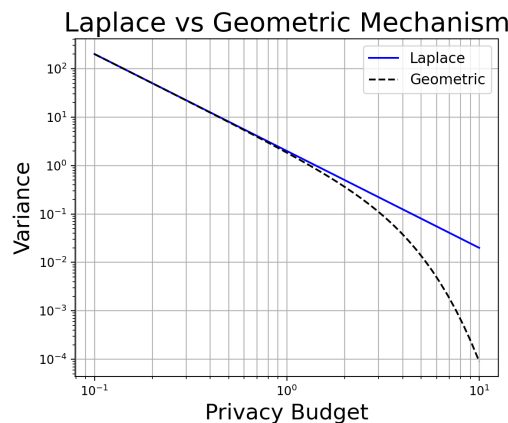


FIGURE 1 – Figure à construire

1.2 Comparaison des variances sur une table

Au transparent 20 du CM2 est présenté un algorithme ϵ -DP permettant de construire un histogramme PVP. L'erreur entre celui-ci et l'original peut être évaluée en calculant la norme-1 entre ces deux histogrammes.

Question 1 :

Pour un ϵ donné proposer une fonction qui construit un histogramme ϵ -DP (sans l'afficher) qui retourne l'erreur entre l'histogramme original et celui-ci ainsi construit.

Question 2 :

Pour un ϵ et un nombre de répétitions n , proposer une fonction qui calcule n fois l'erreur entre les deux histogrammes, puis retourne la moyenne \bar{e} et l'écart-type de cette série d'erreurs σ_e .

Question 3 :

Dans une seconde figure, afficher la courbe des erreurs moyennes \bar{e} pour ϵ qui varie de 0.1 à 10 et griser à chaque fois l'intervalle $[\bar{e} - \sigma_e, \bar{e} + \sigma_e]$. On prendra une échelle logarithmique pour le budget ϵ .

Question 4 :

Implanter le mécanisme géométrique et son application à un histogramme.

Question 5 :

Reprendre le code de la figure précédente pour y ajouter la courbe des erreurs moyennes \bar{e}' de production d'histogrammes avec le mécanisme géométrique pour ϵ qui varie de 0.1 à 10 et griser à chaque fois l'intervalle $[\bar{e}' - \sigma_{e'}, \bar{e}' + \sigma_{e'}]$. On prendra une échelle logarithmique pour le budget ϵ .

2 Critique de la DP pour des requêtes sur certaines tables

Cette partie a pour objectif de publier de manière respectueuse des résultats de requêtes sur une table, en renforçant ce qui a déjà été vu en cours. Elle s'appuie largement sur les sections 1 et 2 de [l'article accessible en ligne](#) :

Wilson, R. J., Zhang, C. Y., Lam, W., Desfontaines, D., Simmons-Marengo, D., & Gipson, B. (2019). Differentially private SQL with bounded user contribution. arXiv preprint arXiv :1909.01917.

On considère la base de donnée D_1 réduite à la table `access_logs` affichée à la Table 1. On reprend les notations de l'article.

alid	uid	browser_agent
1	1	A
2	2	B
3	3	C
4	1	A
5	2	B
6	1	A

TABLE 1 – Table `access_logs`

Question 6 : Trouver une table D_2 de cardinalité inférieure à celle de D_1 telle que $\|D_1 - D_2\| = 1$. Même question avec $\|D_1 - D_2\|_u = 1$.

Question 7 : On considère la requête Q de construction de l'histogramme de fréquence d'utilisation des différents navigateurs qui est formalisée au listing 1 (page 3) de l'article. Quelle serait la réponse d'une telle requête sur la base D_1 ?

Question 8 : Pourquoi la sensibilité ΔQ d'une telle requête vaut-elle 1 ? Pourquoi $\Delta_u Q$ quant-à elle n'est-elle pas bornée ?

Question 9 : Expliquer la phrase "One naive approach is to add Laplace noise of scale $1/\epsilon$ to each count" écrite immédiatement après le listing 1 (page 3). Dans quel contexte cette approche naïve est-elle correcte cependant ?

Question 10 : Pourquoi le remplacement de `COUNT (*)` par `COUNT (DISTINCT uid)` dans le listing 2 (à cheval entre la page 3 et la page 4) permet-il de corriger cette erreur ?

Question 11 : Dans la vraie vie, expliquer pourquoi dans la table `access_logs` un même `uid` pourrait être associé à plusieurs `browser_agent` différents ? Construire une table D'_1 qui étend D_1 et qui possède cette propriété.

Question 12 : Dans le listing 4 (à cheval entre la page 4 et la page 5), expliquer les instructions :

1. TABLESAMPLE RESERVOIR (Cu **ROWS** PARTITION **BY** uid) ;
2. **COUNT** (**DISTINCT** uid) + Laplace (Cu/epsilon).